



# Query-Oriented Summarization of RDF Graphs

Šejla Čebirić, François Goasdoué, Ioana Manolescu

## ► To cite this version:

Šejla Čebirić, François Goasdoué, Ioana Manolescu. Query-Oriented Summarization of RDF Graphs. Data Science - 30th British International Conference on Databases, BICOD 2015, Edinburgh, UK, July 6-8, 2015, Proceedings, Jul 2015, Edinburgh, United Kingdom. pp.87–91. hal-01176277

**HAL Id: hal-01176277**

**<https://hal.inria.fr/hal-01176277>**

Submitted on 15 Jul 2015

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Query-Oriented Summarization of RDF Graphs

Šejla Čebirić<sup>1</sup>, François Goasdoué<sup>2,1</sup>, and Ioana Manolescu<sup>1</sup>

<sup>1</sup> INRIA and U. Paris-Sud, France {sejla.cebirc,ioana.manolescu}@inria.fr

<sup>2</sup> U. Rennes 1, Lannion, France fg@irisa.fr

**Abstract.** The Resource Description Framework (RDF) is the W3C’s graph data model for Semantic Web applications. We study the problem of RDF graph summarization: given an input RDF graph  $G$ , find an RDF graph  $S_G$  which summarizes  $G$  as accurately as possible, while being possibly orders of magnitude smaller than the original graph. Our approach is *query-oriented*, i.e., querying a summary of a graph should reflect whether the query has some answers against this graph. The summaries are aimed as a help for query formulation and optimization. We introduce two summaries: a *baseline* which is compact and simple and satisfies certain accuracy and representativeness properties, but may oversimplify the RDF graph, and a *refined* one which trades some of these properties for more accuracy in representing the structure.

## 1 Introduction

The Resource Description Framework (RDF) is a graph-based data model promoted by the W3C as the standard for Semantic Web applications; SPARQL is the W3C’s standard language for querying RDF data.

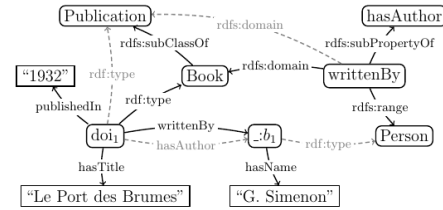
RDF graphs are often *large* and *varied*, produced in a variety of contexts, e.g., scientific applications, social or online media, government data etc. They are *heterogeneous*, i.e., resources described in an RDF graph may have very different sets of properties. An RDF resource may have: no types, one or several types (which may or may not be related to each other). *RDF Schema* (RDFS) information may optionally be attached to an RDF graph, to enhance the description of its resources. Such statements also entail that in an RDF graph, some data is **implicit**. According to the W3C RDF and SPARQL specification, **the semantics of an RDF graph comprises both its explicit and its implicit data**; in particular, SPARQL query answers must be computed *reflecting both its explicit and implicit data*, even if the latter is not physically stored.

In this work, we study the problem of *RDF graph summarization*, that is: given an input RDF graph  $G$ , find an RDF graph  $S_G$  which *summarizes  $G$  as accurately as possible, while being possibly orders of magnitude smaller* than the original graph. Such a summary can be used in a variety of contexts: to help an RDF application designer get acquainted with a new dataset, as a first-level user interface, or as a support for query optimization as traditionally the case in semi-structured graph data management [3] etc. Our approach is *query-oriented*, i.e., querying a summary of a graph should reflect whether the query has some answers against this graph. The properties our summaries aim at are related to query processing, in particular enabling static analysis, query formulation

and optimization (i.e., deciding if a query is empty or finding a simpler way to formulate a query). While semi-structured data summarization has been studied before, our work is the first focused on *partially explicit, partially implicit* RDF graphs. Our ongoing technical report [5] provides proofs for the results presented here and a discussion of related work.

## 2 RDF Graphs and Summary Requirements

**RDF graphs and queries** An *RDF graph* (or *graph*) is a set of *triples* of the form  $s \ p \ o$ , stating that the *subject*  $s$  has the *property*  $p$ , and the value of that property is the *object*  $o$ . Triples are formed using uniform resource identifiers (URIs), typed or untyped literals (constants), and *blank nodes* (unknown URIs or literals) corresponding to incomplete information. We use  $s$ ,  $p$ , and  $o$  in triples as placeholders. Literals are shown as strings between quotes, e.g., “string”.



**Fig. 1.** Sample RDF graph

author (a blank node  $_:b_1$  related to the author name), title and publication date.

$$G = \{doi_1 \text{ rdf:type Book}, doi_1 \text{ writtenBy } _:b_1, doi_1 \text{ publishedIn "1932"}, \\ doi_1 \text{ hasTitle "Port des Brumes"}, _:b_1 \text{ hasName "G. Simonon"}\}$$

**RDF Schema** triples allow enhancing the descriptions in RDF graphs by declaring *deductive constraints* between the graph classes and properties, namely: *subClassOf*, *subPropertyOf*, *domain* and *range*, where the latter two denote the first and second attribute of a property, respectively. Consequently, an RDF graph may have **implicit triples** even though they do not exist explicitly. For instance, assume the RDF graph  $G$  above is extended with the following constraints:

- books are publications: **Book** *rdfs:subClassOf* **Publication**
- writing something means being an author:  
**writtenBy** *rdfs:subPropertyOf* **hasAuthor**
- writtenBy is a relation between books and people:  
**writtenBy** *rdfs:domain* **Book** and **writtenBy** *rdfs:range* **Person**

The resulting graph is depicted in Fig. 1. Its implicit triples are those represented by dashed-line edges. Adding all the implicit triples to an RDF graph  $G$  leads to its *saturation*  $G^\infty$ , which is the RDF graph stating the semantics of  $G$ .

In this work, we consider *conjunctive* SPARQL queries, a.k.a. Basic Graph Pattern (BGP) queries. The evaluation of a query  $q$  against an RDF graph  $G$  based on  $G$ ’s explicit triples may lead to an incomplete answer; the complete answer is obtained by evaluating  $q$  against  $G^\infty$ . E.g., consider:

$q(x_3) :- x_1 \text{ hasAuthor } x_2, x_2 \text{ hasName } x_3, x_1 \text{ hasTitle "Le Port des Brumes"}$

Its answer against the graph in Fig. 1 is  $q(\mathbf{G}^\infty) = \{\langle \text{“G. Simenon”} \rangle\}$ . Note that evaluating  $q$  against  $\mathbf{G}$  leads to an empty answer.

**Summary requirements** We assume that *the summary  $\mathbf{S}_\mathbf{G}$  of an RDF graph  $\mathbf{G}$  is an RDF graph itself*. Further, we require summaries to satisfy the following conditions: (i) The saturation of the summary of an RDF graph  $\mathbf{G}$  must be the same as the summary of its saturation  $\mathbf{G}^\infty$ , since the semantics of an RDF graph is its saturation; (ii) The summary should be (if possible, much) smaller than the RDF graph; (iii) The summary should be *representative*: queries with results on  $\mathbf{G}$  should also have results on the summary; (iv) The summary should be *accurate*: queries with results on the summary should reflect that such data existed indeed in the graph. To formalize these, let  $\mathcal{Q}$  be a SPARQL dialect.

**Definition 1.** (*Query-Based Representativeness*)  $\mathbf{S}_\mathbf{G}$  is  $\mathcal{Q}$ -representative of  $\mathbf{G}$  if and only if for any query  $q \in \mathcal{Q}$  such that  $q(\mathbf{G}^\infty) \neq \emptyset$ , we have  $q(\mathbf{S}_\mathbf{G}^\infty) \neq \emptyset$ .

Note that *several graphs may have the same summary*, since a summary loses *some* of the information from the original graph. If two RDF graphs differ only with respect to such information, they have the same summary. We term *inverse set* of  $\mathbf{S}_\mathbf{G}$ , the set of all RDF graphs whose summary is  $\mathbf{S}_\mathbf{G}$ . This leads to the accuracy criterion, with respect to *any graph a summary may correspond to*:

**Definition 2.** (*Query-Based Accuracy*) Let  $\mathbf{S}_\mathbf{G}$  be a summary, and  $\mathcal{G}$  the inverse set of  $\mathbf{S}_\mathbf{G}$ . The summary  $\mathbf{S}_\mathbf{G}$  is  $\mathcal{Q}$ -accurate if for any query  $q \in \mathcal{Q}$  such that  $q(\mathbf{S}_\mathbf{G}^\infty) \neq \emptyset$ , there exists  $\mathbf{G} \in \mathcal{G}$  such that  $q(\mathbf{G}^\infty) \neq \emptyset$ .

For compactness, the (voluminous) set of literals, along with subject and object URIs for non-**type** triples from  $\mathbf{G}$  should not appear in the summary. However, given that property URIs are often specified in SPARQL queries [1], and that typically there are far less distinct property URIs than the subject or object URIs [4], property URIs should be preserved by the summary. This leads us to considering the following SPARQL dialect:

**Definition 3.** (*RBGP queries*) A relational basic graph pattern query (RBGP) is a conjunctive SPARQL query whose body has: (i) URIs in all the property positions, (ii) a URI in the object position of every **type** triple, and (iii) variables in any other positions.

We define *RBGP representativeness* and *RBGP accuracy* by instantiating  $\mathcal{Q}$  in Definition 1 and Definition 2, respectively, to RBGP queries.

### 3 RDF Summaries

We assume a function `newURI()` returning a fresh URI on each call. We call *data property* any property  $p$  in  $\mathbf{G}$  different from **type**. Further, for any data property  $p$ , the *property source* of  $p$ , denoted  $\mathcal{S}(p)$ , is a URI set using `newURI()`, and similarly, the *property target* of  $p$ , denoted  $\mathcal{T}(p)$ , is a URI set using `newURI()`.

We introduce our summaries below; examples are delegated to [5] and can also be found at <https://team.inria.fr/oak/projects/rdfsummary/>.

**Definition 4.** (*Baseline Summary*) Given an RDF graph  $\mathbf{G}$ , the baseline summary of  $\mathbf{G}$  is an RDF graph  $\mathbf{B}_\mathbf{G}$  such that:

**Schema**  $B_G$  has the same schema triples as  $G$ .

**DNT** (Data triples of  $B_G$  whose property is not **type**) Let  $p, p_1, p_2$  be some data properties from  $G$ .

**DNT1** The triple  $S(p) \ p \ T(p)$  belongs to  $B_G$ ;

**DNT2** if  $s \ p_1 \ o_1, s \ p_2 \ o_2 \in G$ , then  $S(p_1) = S(p_2)$ ;

**DNT3** if  $s_1 \ p_1 \ o, s_2 \ p_2 \ o \in G$ , then  $T(p_1) = T(p_2)$ ;

**DNT4** if  $s \ p_1 \ o_1, o_1 \ p_2 \ o_2 \in G$ , then  $T(p_1) = S(p_2)$ ;

**DT** (Data triples of  $B_G$  whose property is **type**)

**DT1** If  $s \ p \ o, s \ \text{type } c$  are in  $G$ , then  $S(p) \ \text{type } c$  is in  $B_G$ ;

**DT2** if  $s \ p \ o, o \ \text{type } c$  are in  $G$ , then  $T(p) \ \text{type } c$  is in  $B_G$ ;

**DT3** Let  $n_{all}$  be set to `newURI()`. If  $s \ \text{type } c \in G$ , and  $\nexists s \ p \ o \in G$ , then  $n_{all} \ \text{type } c \in B_G$ .

**Refined summary** The baseline summary may unify property source and target URIs quite aggressively. For instance, if a *store* and a *person* both have a *zipcode*, they will lead to the same baseline URI, even though they are very different things. To mitigate this issue, we designed a second flavor of summary of an RDF graph  $G$ , termed *refined* and denoted  $R_G$ . For space reasons, the definition is delegated to [5]. Intuitively, the difference between the baseline and the refined summary is that the latter fuses data property source and/or target URIs *only if one resource in  $G$  that leads to their unification has no type at all*.

**Summary properties** Both summaries meet our requirements (i), (iii) and (iv) as follows. We say two summary graphs are *equivalent*, denoted  $\equiv$ , iff they are identical up to a bijection between their sets of URIs. The summaries *commute with saturation*, i.e.,  $(S_G)^\infty \equiv S_{G^\infty}$ , and are RBGP accurate. The  $B_G$  is fully RBGP representative, and the  $R_G$  is representative of *RBGPs having no more than one type triple with the same subject*. This follows from a *graph homomorphism* from  $G^\infty$  to  $(S_G)^\infty$  [5]. Observe that  $S_G$  is not a core of  $G$ , since we cannot guarantee a homomorphism from  $S_G$  to  $G$  ( $S_G$  may comprise false positives).

The size of the baseline summary is bounded by the size of  $G$ 's schema plus the number of data properties and class assertions from  $G$ . It can be built in  $O(|G|^2)$  time. Computing the refined summary has  $O(|G|^5)$  complexity, requiring an efficient underlying system e.g., based on triple partitioning and indexing or a distributed processing platform such as [2]. An upper bound for its size is the number of classes in  $G \times$  the number of distinct data properties in  $G$ .

**Acknowledgments** This work has been partially funded by the projects Datalyse "Investissement d'Avenir" and ODIN "DGA RAPID".

## References

- [1] M. Arias, J. D. Fernández, M. A. Martínez-Prieto, and P. de la Fuente. An empirical study of real-world SPARQL queries. *CoRR*, abs/1103.5043, 2011.
- [2] F. Goasdoué, Z. Kaoudi, I. Manolescu, J.-A. Quiané-Ruiz, and S. Zampetakis. CliqueSquare: Flat Plans for Massively Parallel RDF Queries. In *ICDE*, 2015.
- [3] R. Goldman and J. Widom. Dataguides: Enabling query formulation and optimization in semistructured databases. In *VLDB*, 1997.
- [4] Statistics on The Billion Triple Challenge Dataset. [gromgull.net/blog/2010/09/btc2010-basic-stats](http://gromgull.net/blog/2010/09/btc2010-basic-stats), 2010.
- [5] Technical report, 2015.